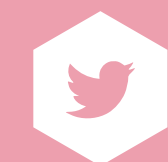


# DOCKER AND PYTHON

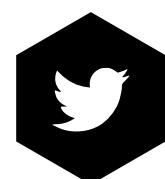
**Making them play nicely and securely for Data Science and Machine Learning**

TANIA ALLARD, PHD

Sr. Developer Advocate @Microsoft.



ixek | <https://bit.ly/europython-ml-docker>



**@ixek**



**@trallard**



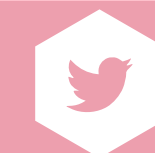
**trallard.dev**

# THESE SLIDES

<https://bit.ly/europython-ml-docker>

# WHAT YOU'LL LEARN TODAY

- Why using Docker?
- Docker for Data Science and Machine Learning
- Security and performance
- Do not reinvent the wheel, automate
- Tips and trick to use Docker

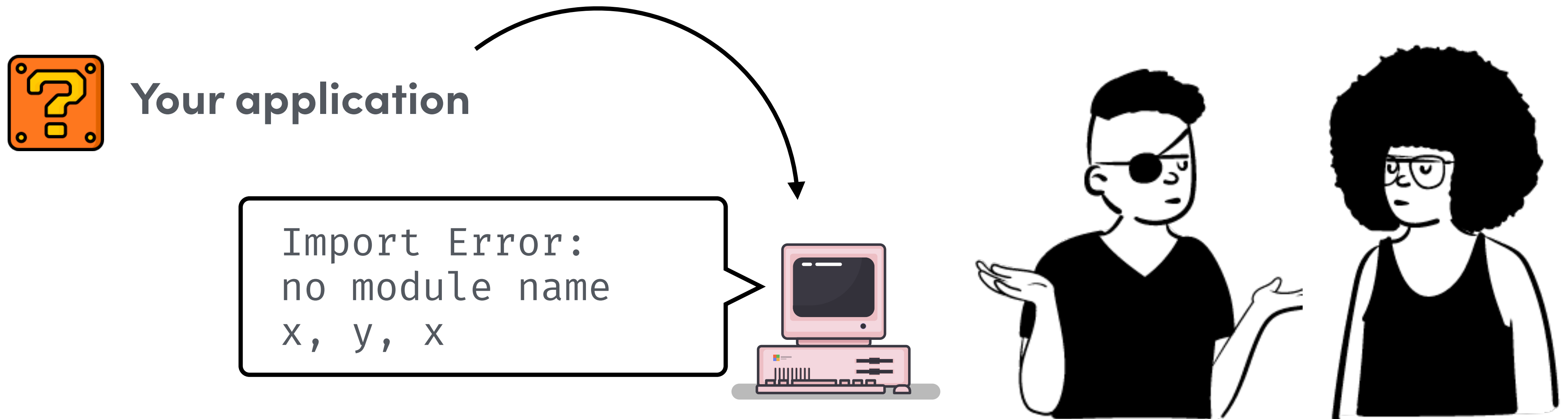




**WHY**

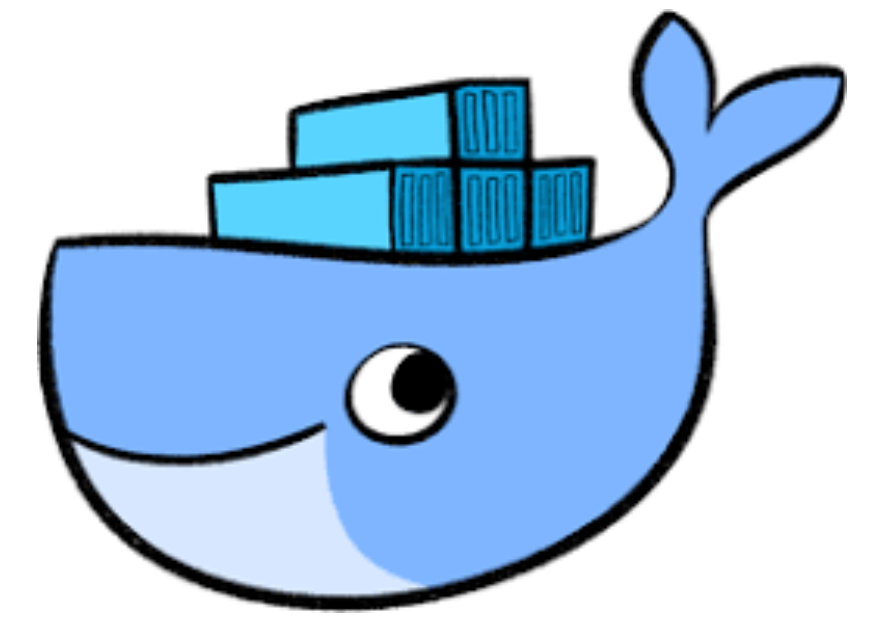
**DOCKER?**

# DEV LIFE WITHOUT DOCKER OR CONTAINERS



**How are your users or colleagues meant to know what dependencies they need?**

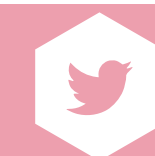
# WHAT IS DOCKER?



A tool that helps you to create, deploy and run your applications or projects **by using containers.**

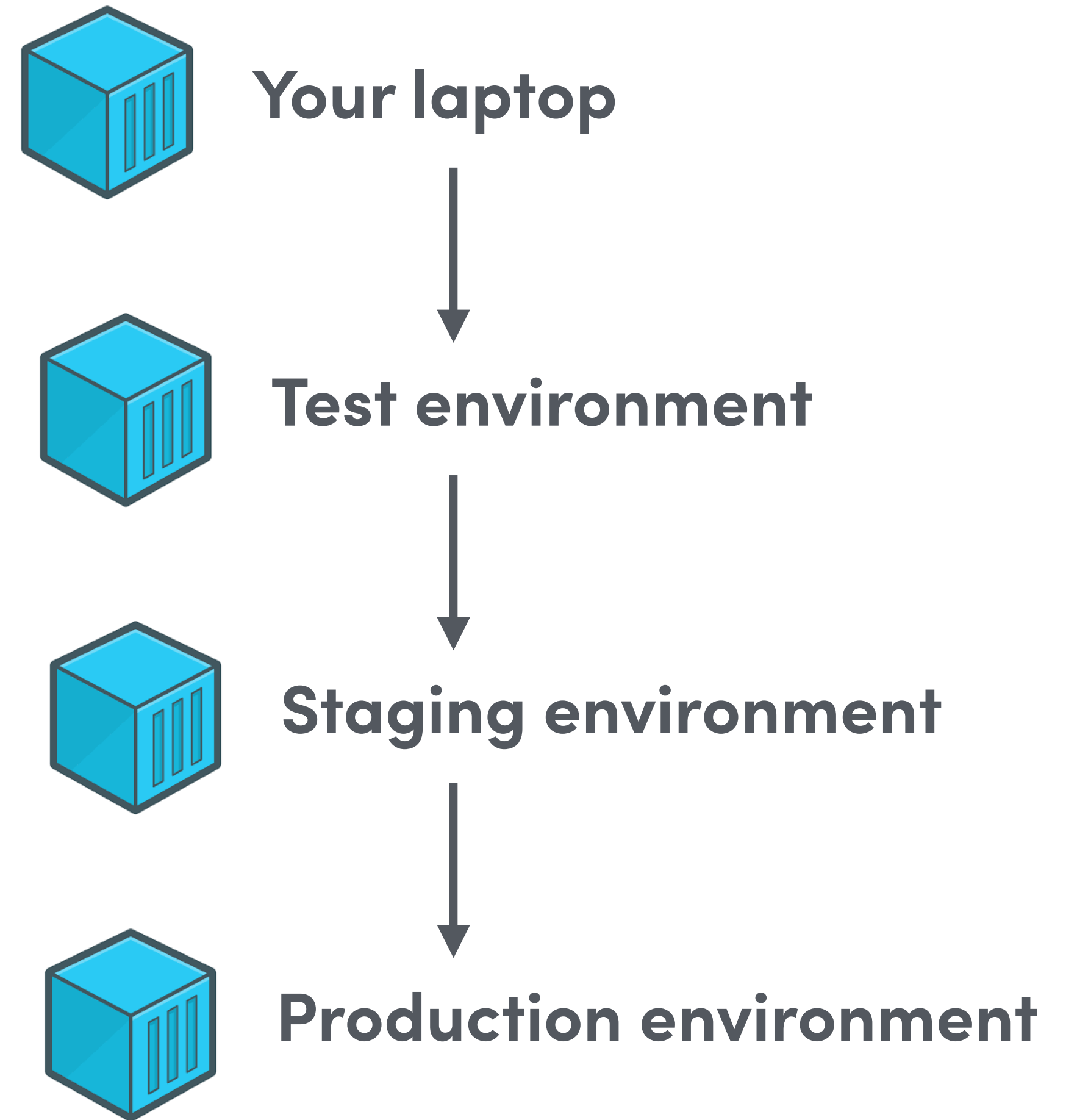


This is a container



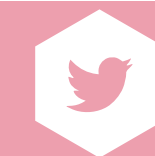
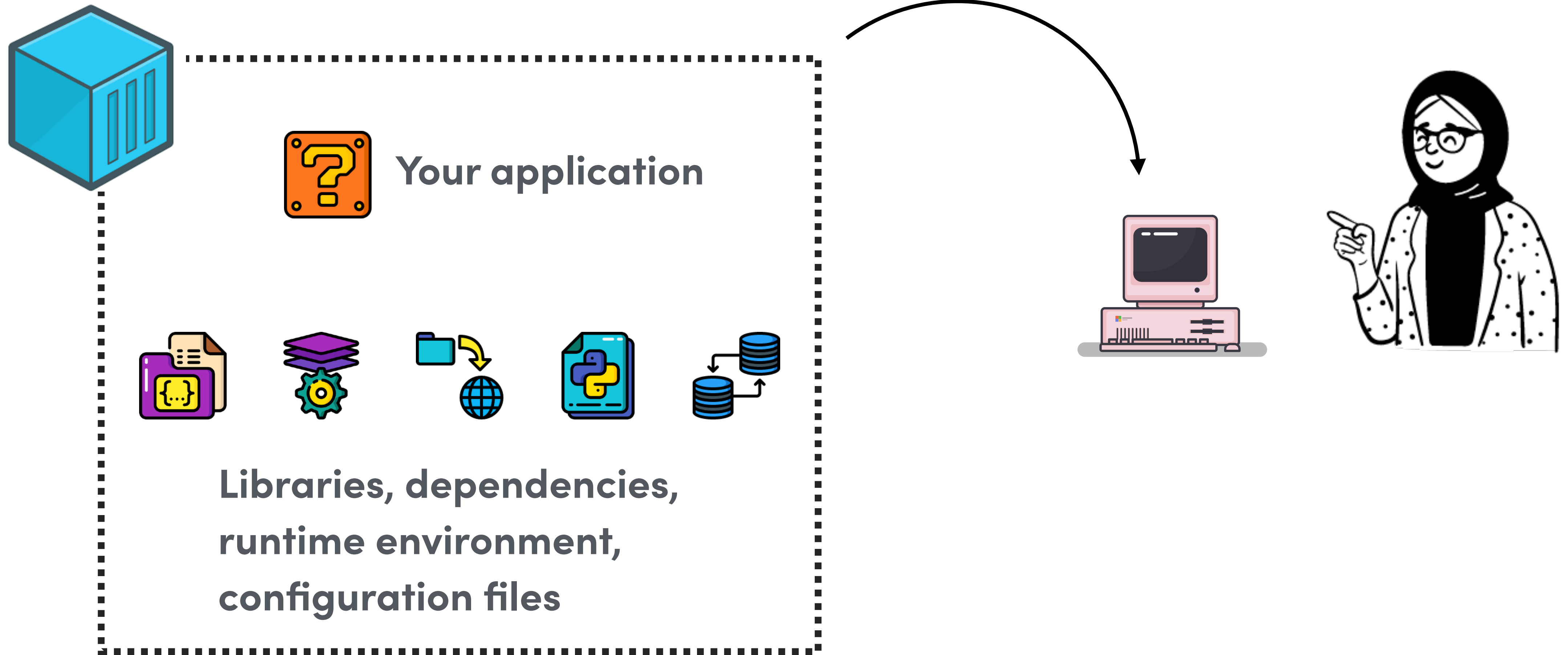
# HOW DO CONTAINERS HELP ME?

They provide a solution to the problem of how to get software to run reliably when **moved from one computing environment to another**



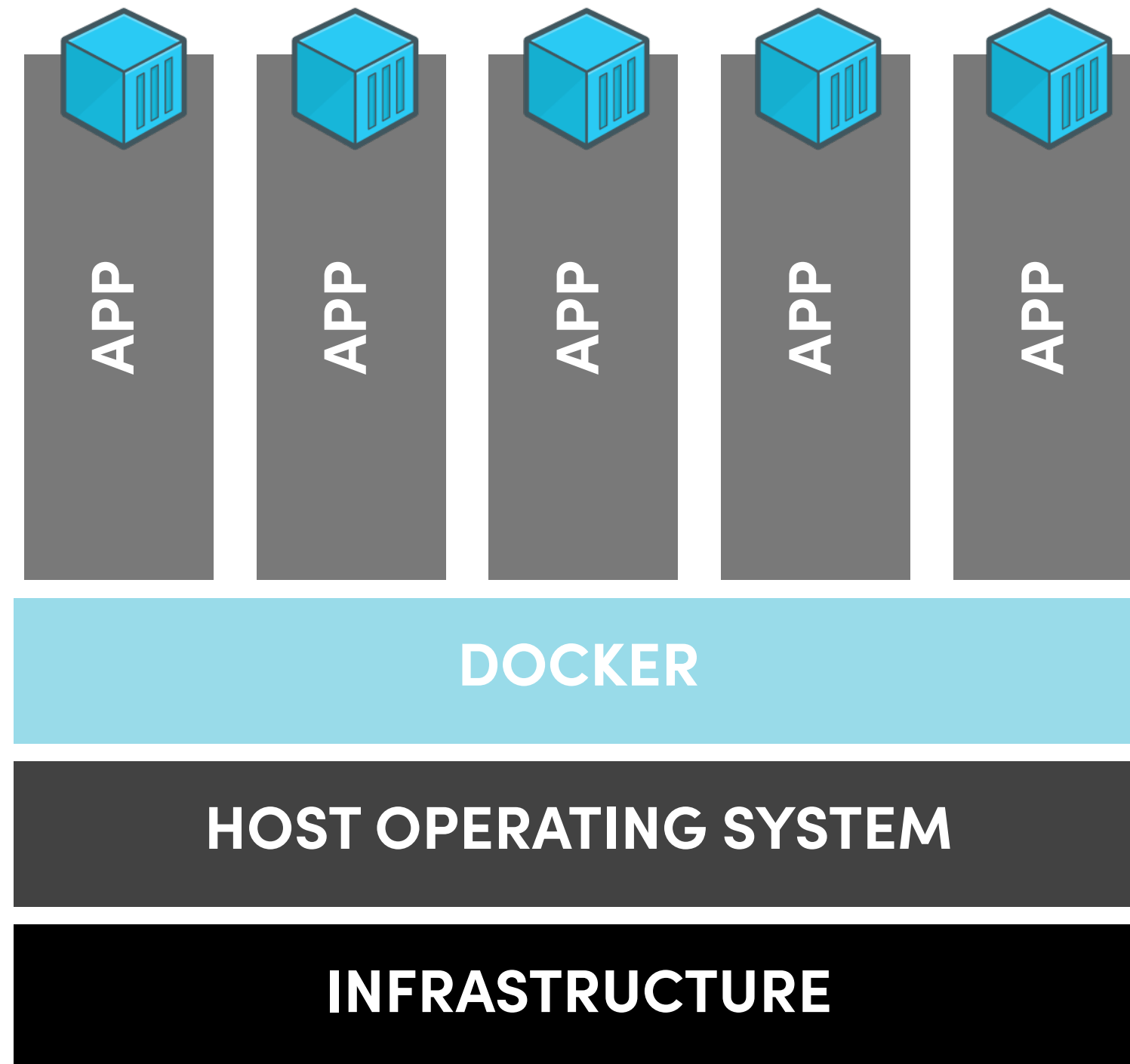


# DEV LIFE WITH CONTAINERS

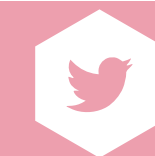


# THAT SOUNDS A LOT LIKE A VIRTUAL MACHINE

Each app is  
containerised

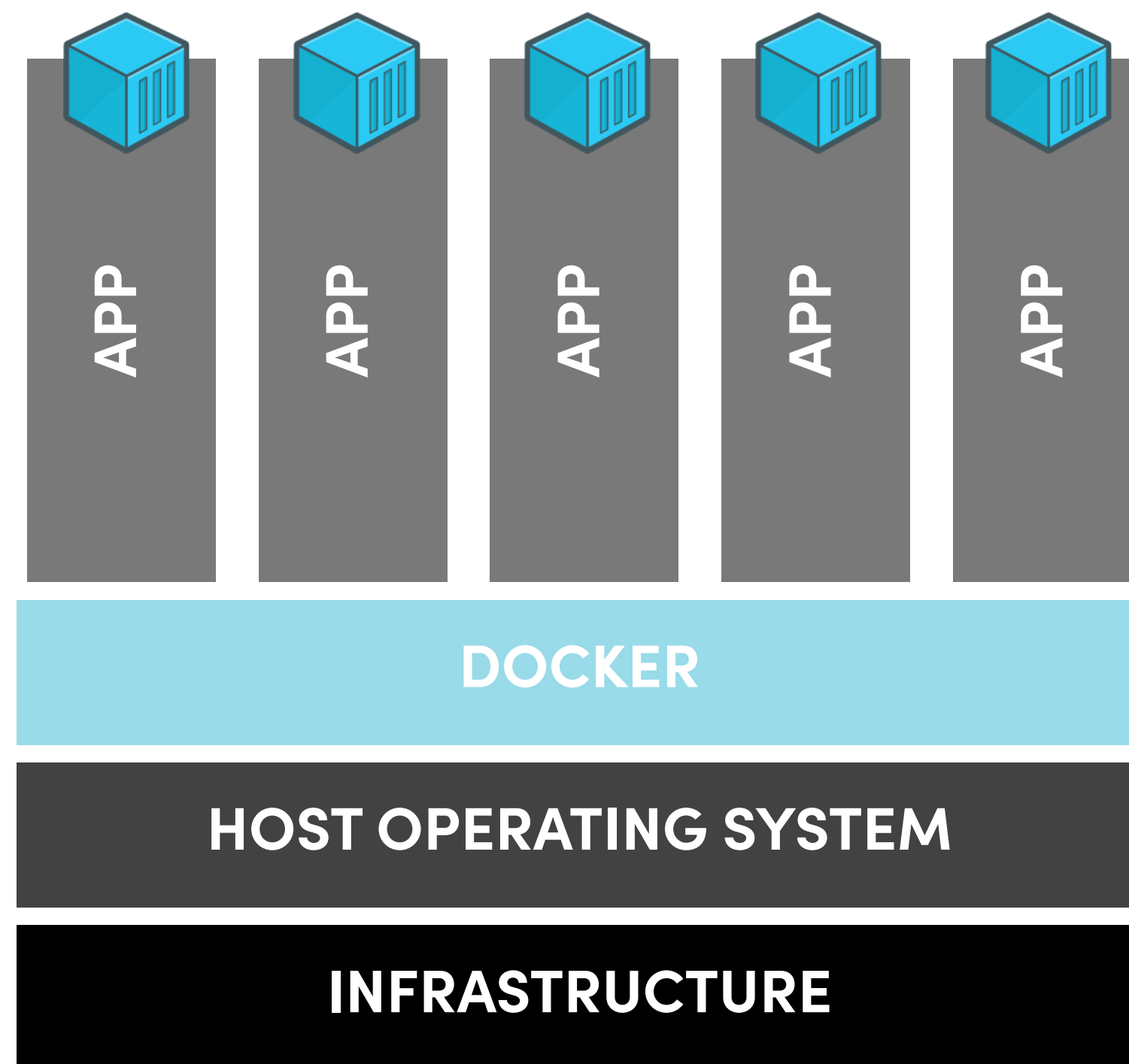


At the app level:  
Each runs as an isolated process



# THAT SOUNDS A LOT LIKE A VIRTUAL MACHINE

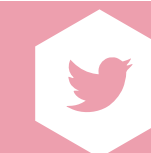
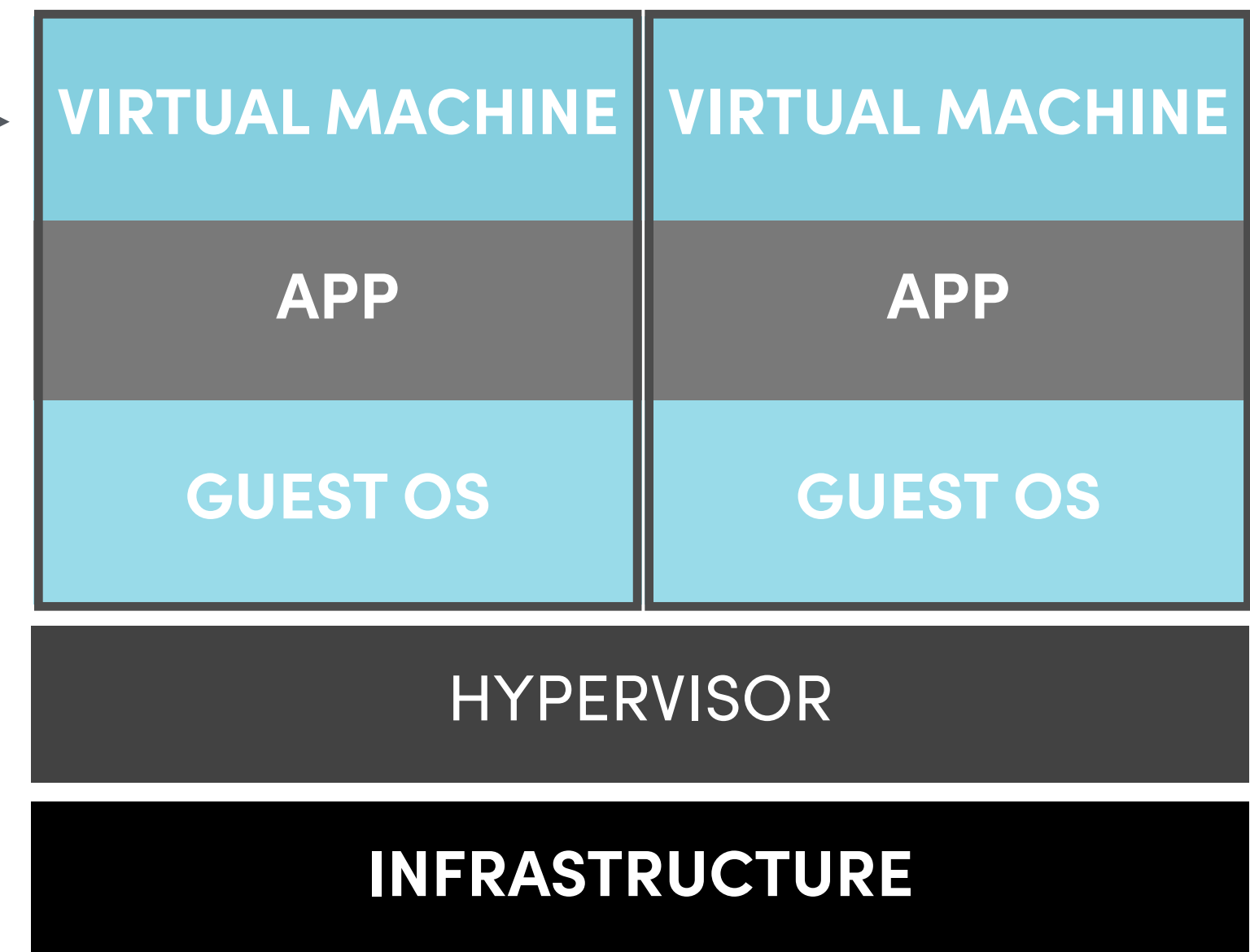
## CONTAINERS



## VIRTUAL MACHINE

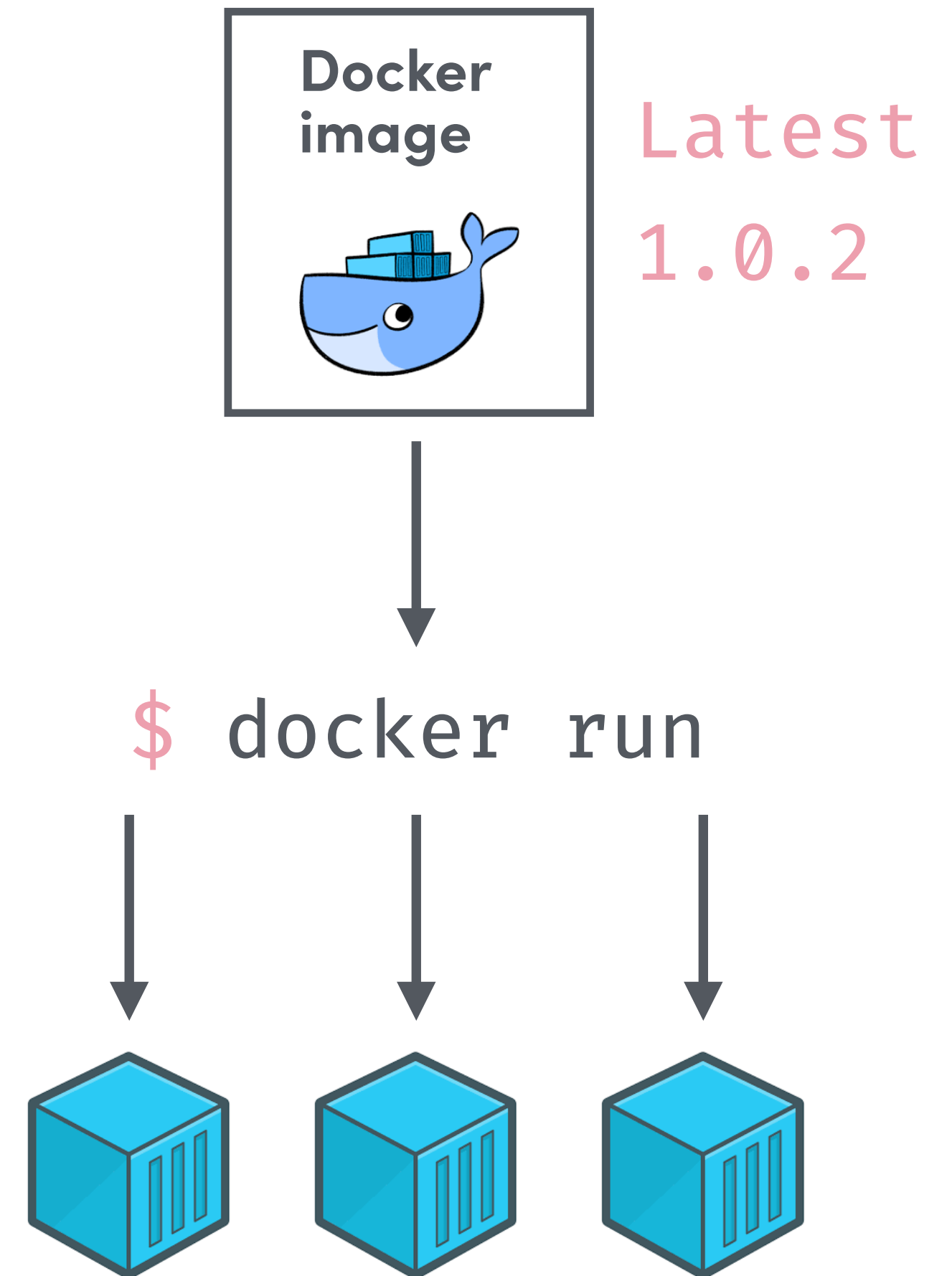
At the hardware level

Full OS + app + binaries + libraries



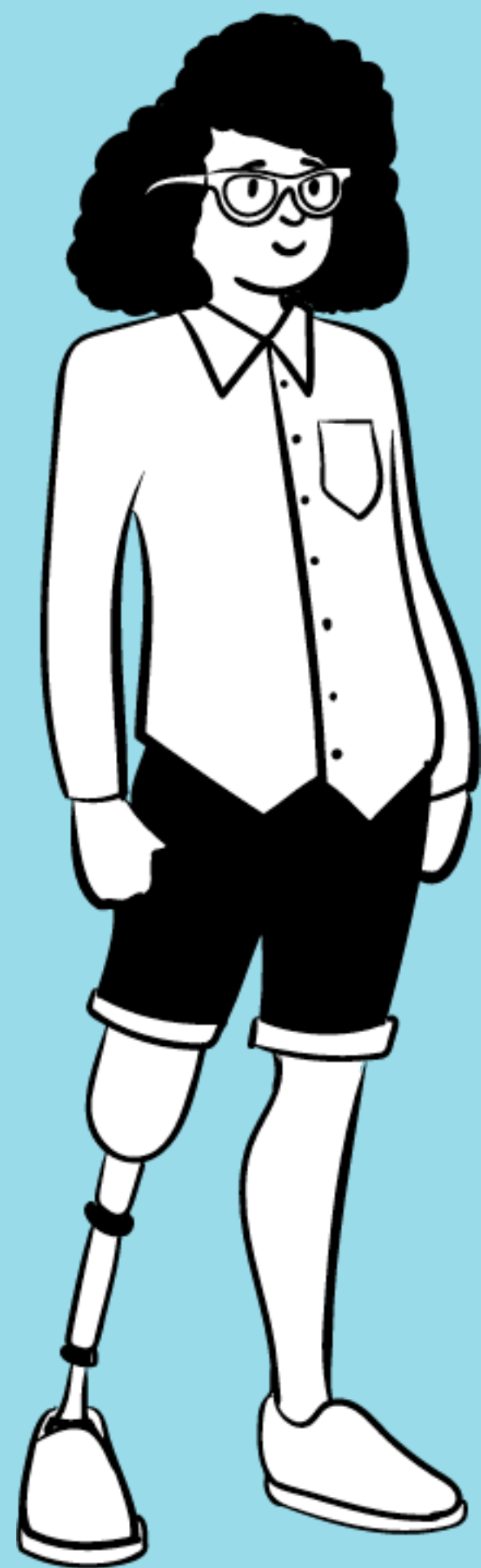
# IMAGE VS CONTAINER

- **Image**: archive with all the data needed to run the app
- When you run an image it creates a **container**



# COMMON PAIN POINTS IN DS AND ML

- Complex setups / dependencies
- Reliance on data / databases
- Fast evolving projects (iterative R&D process)
- Docker is complex and can take a lot of time to upskill
- Are containers secure enough for my data / model /algorithm?



# DOCKER FOR DATA SCIENCE AND MACHINE LEARNING

# HOW IS IT DIFFERENT FROM WEB APPS FOR EXAMPLE?



## Drawing a line to the scope of Python packaging

■ Packaging



Tzu-ping Chung uranusjr

Feb '19

Another topic in the Big Picture thread I found interesting 😊

twitter.com



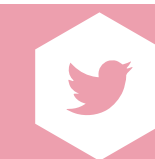
**Sylvain Corlay (SylvainCorlay)**

@dstufft @pwang @WillingCarol @brettsky @vorpalsmith @zooba @uranusjr @acanthamoeba @ncoghlan\_dev @kushaldas Although there is a continuum of things between "we need BLAS", and "we need R". Where is the limit?

It feels to me that with language-specific packaging tools, you soon find a cliff somewhere between those two usecases.

6:46 PM - 11 Feb 2019

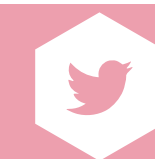
<https://twitter.com/dstufft/status/1095164069802397696>



ixek | <https://bit.ly/euopython-ml-docker>

# HOW IS IT DIFFERENT FROM WEB APPS FOR EXAMPLE?


- Not every deliverable is an app
- Not every deliverable is a model either
- Heavily relies on data
- Mixture of wheels and compiled packages
- Security access levels - for data and software
- Mixture of stakeholders: data scientists, software engineers, ML engineers



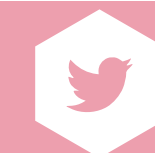


# BUILDING DOCKER IMAGES

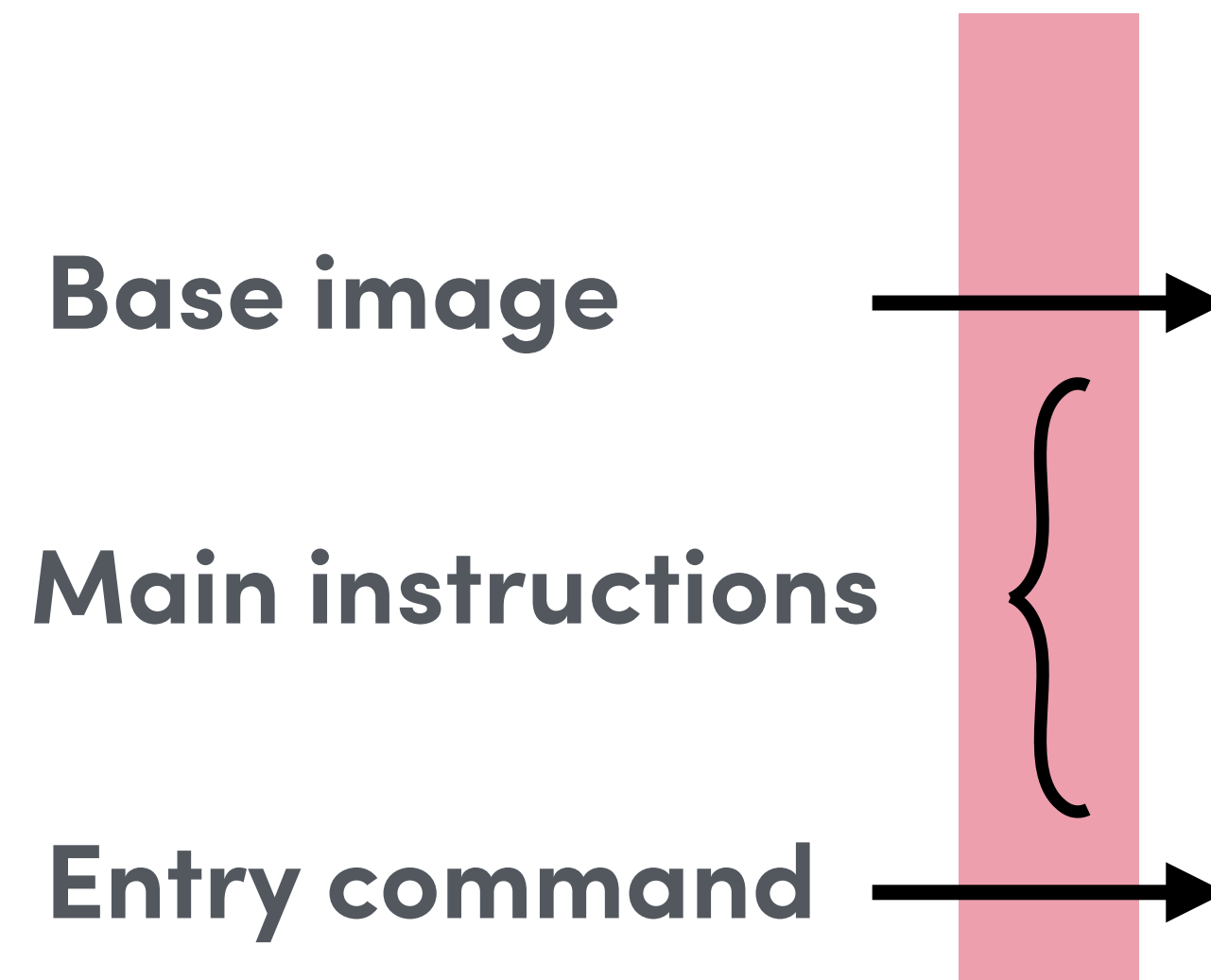
**Dockerfiles** are used to create Docker images by providing a set of instructions to install software, configure your image or copy files



```
docker-ds - Dockerfile
1 # word of caution - this is a bad example
2 FROM python:3
3
4 COPY yourscript.py /
5
6 RUN pip install flask
7
8 CMD [ "python", "./yourscript.py" ]
```

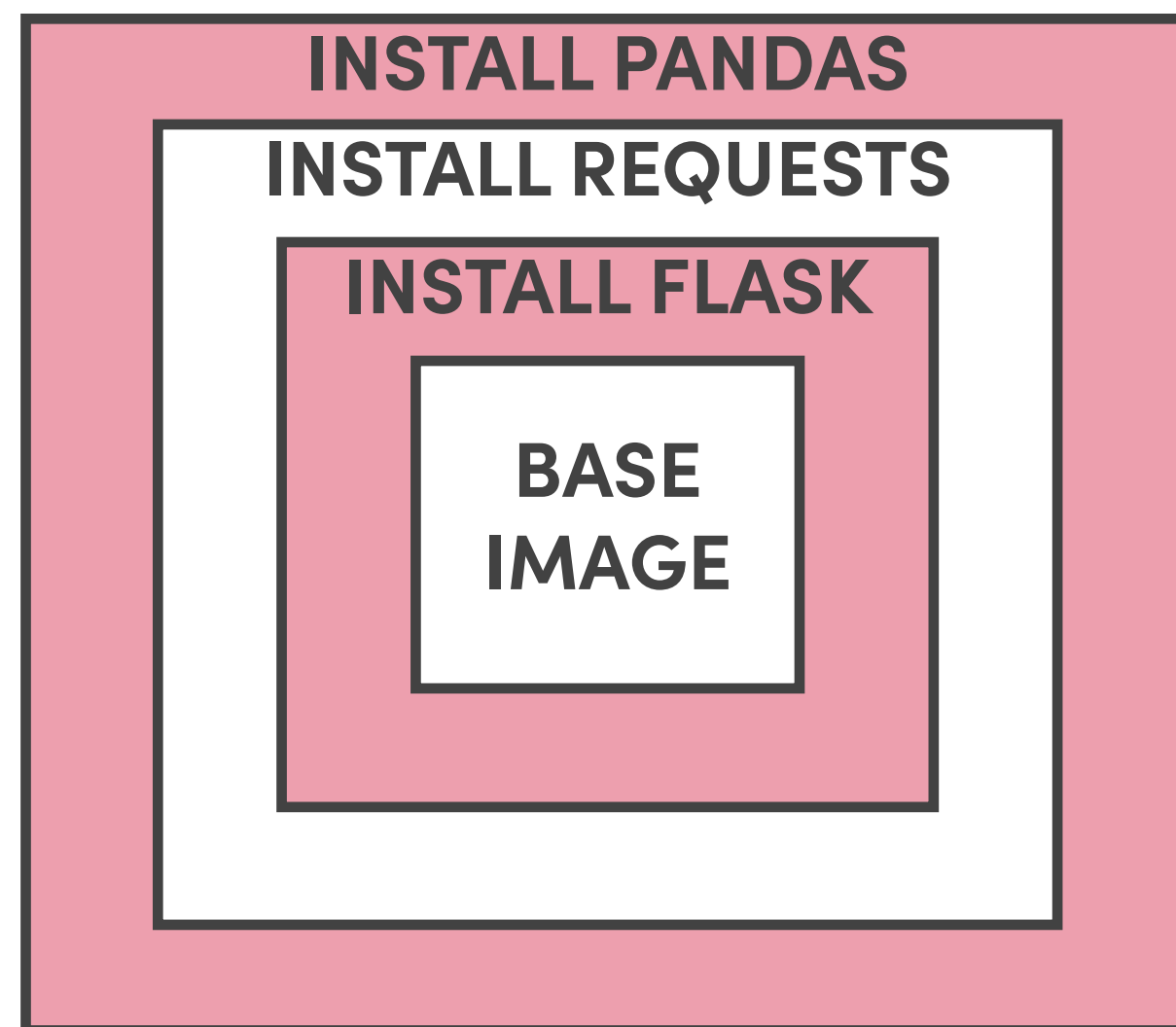


# DISSECTING DOCKER IMAGES



```
docker-ds - Dockerfile
1 # word of caution - this is a bad example
2 FROM python:3
3
4 COPY yourscript.py /
5
6 RUN pip install flask
7
8 CMD [ "python", "./yourscript.py" ]
```

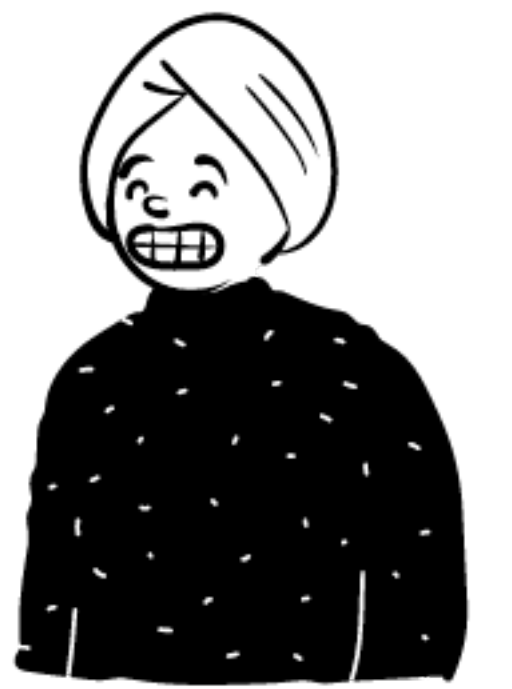
# DISSECTING DOCKER IMAGES



```
docker-ds - Dockerfile
1  # word of caution - this is a bad example
2  FROM python:3
3
4  COPY yourscrip.py /
5
6  RUN pip install flask
7  RUN pip install requests
8  RUN pip install pandas
9
10 CMD [ "python", "./yourscrip.py" ]
```

Each instruction creates  
A layer  
(like an onion)

# CHOOSING THE BEST BASE IMAGE

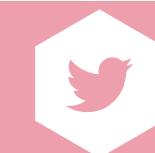


If building from scratch use the **official Python images**

1	REPOSITORY	TAG	SIZE
2	python	3.7.7-alpine	96MB
3	python	3.7.7-slim-stretch	155MB
4	python	3.7.7-stretch	942MB
5	python	3.7.7-slim-buster	179MB
6	python	3.7.7-buster	919MB
7	python	3.8.2-slim-buster	194MB
8	python	3.8.2-buster	934MB

[https://hub.docker.com/\\_/python](https://hub.docker.com/_/python)

<https://github.com/docker-library/docs/tree/master/python>

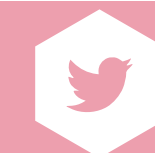
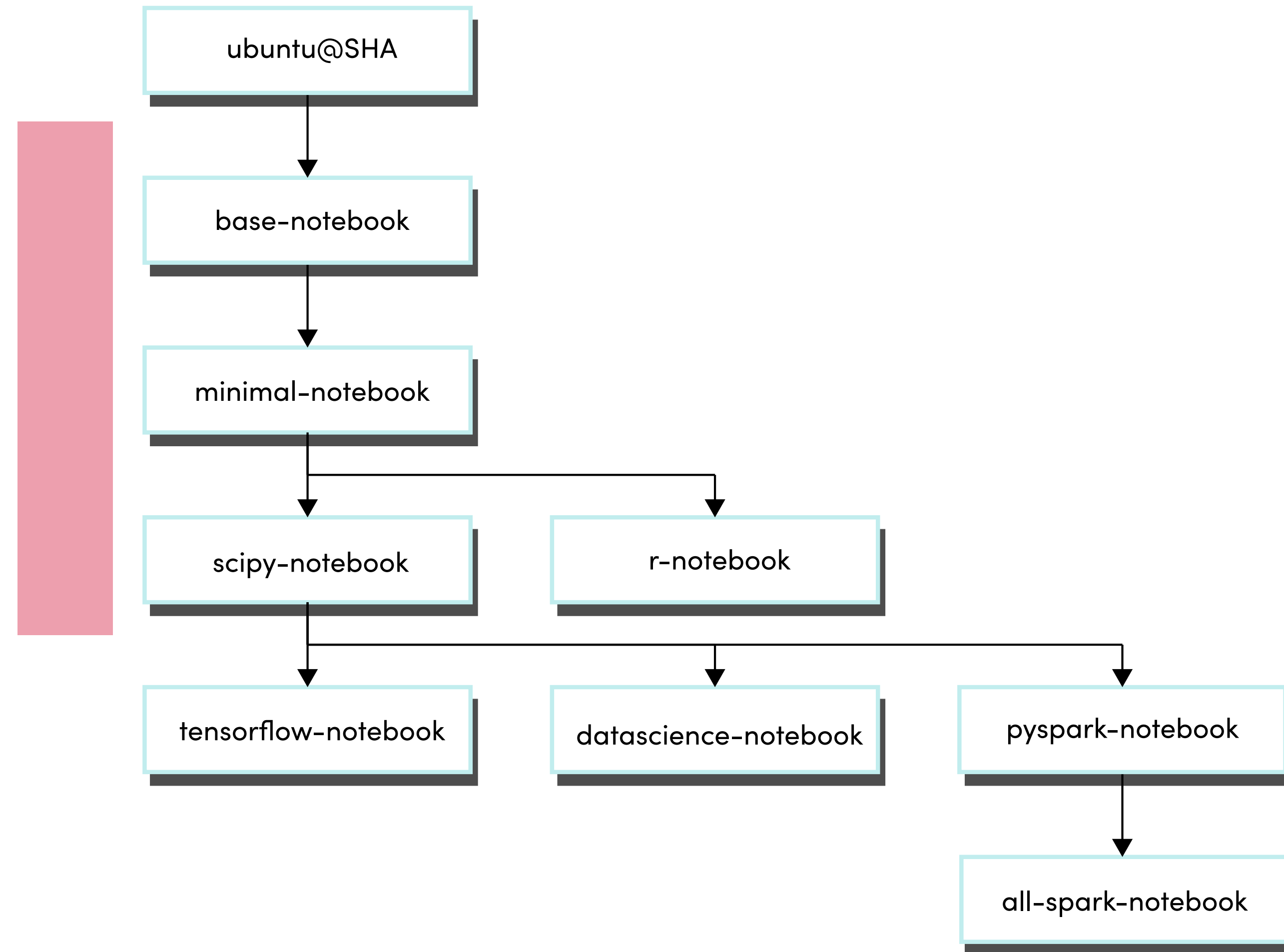


# THE JUPYTER DOCKER STACK

Need Conda, notebooks and scientific Python ecosystem?

Try **Jupyter Docker stacks**

<https://jupyter-docker-stacks.readthedocs.io/>



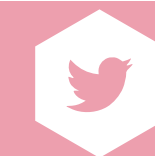
# BEST PRACTICES



- Always know what you are expecting
- Provide context with LABELS
- Split complex RUN statements and sort them
- Prefer COPY to add files

```
docker-ds - Dockerfile
1 # Always use a concrete tag (avoid LATEST)
2 FROM jupyter/base-notebook:6.0.3
3
4 # Add metadata
5 LABEL maintainer="Tania Allard"
6 LABEL securitytxt="https://www.example.com/.well-known/security.txt"
7
8 # Use pinned versions always
9 RUN conda install --quiet --yes \
10     'pandas=1.0.3' \
11     'dask=2.14.*' \
12     && \
13     # do not forget to clean - reduce image
14     conda clean --all -f -y
15
16 # separate instructions per scope
17 RUN mkdir data-sci-demo
18
19 COPY ./your-project data-sci-demo/
20
21
```

[https://docs.docker.com/develop/develop-images/dockerfile\\_best-practices/](https://docs.docker.com/develop/develop-images/dockerfile_best-practices/)



ixek | <https://bit.ly/europython-ml-docker>



# SPEED UP YOUR BUILD

- Leverage build cache
- Install only necessary packages

```
docker-ds - requirements.txt
1 pandas=1.0.3
2 dask=2.14.*
```

```
docker-ds - Dockerfile
1 # Always use a concrete tag (avoid LATEST)
2 FROM jupyter/base-notebook:6.0.3
3
4 # Add metadata
5 LABEL maintainer="Tania Allard"
6 LABEL securitytxt="https://www.example.com/.well-known/security.txt"
7
8 # Leveraging build cache
9 COPY ./requirements.txt /tmp/
10
11 RUN conda install --quiet --yes --file /tmp/requirements.txt &&\
12     # do not forget to clean - reduce image
13     conda clean --all -f -y
14
15 # Separate instructions per scope
16 RUN mkdir data-sci-demo
17
18 COPY ./your-project data-sci-demo/
```

[https://docs.docker.com/develop/develop-images/dockerfile\\_best-practices/](https://docs.docker.com/develop/develop-images/dockerfile_best-practices/)

# SPEED UP YOUR BUILD AND PROOF

- Leverage build cache
- Install only necessary packages
- Explicitly ignore files

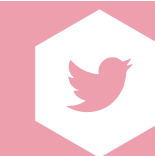
docker-ds - .dockerignore

```
1 # Documentation
2 Readme.md
3
4 # Never add data
5 ./yourproject/data/
6
7 # Secrets
8 appsettings.json
9 .env
10 supersecretkeys.json
```

docker-ds - Dockerfile

```
1 # Always use a concrete tag (avoid LATEST)
2 FROM jupyter/base-notebook:6.0.3
3
4 # Add metadata
5 LABEL maintainer="Tania Allard"
6 LABEL securitytxt="https://www.example.com/.well-known/security.txt"
7
8 # Leveraging build cache
9 COPY ./requirements.txt /tmp/
10
11 RUN conda install --quiet --yes --file /tmp/requirements.txt &&\
12     # do not forget to clean - reduce image
13     conda clean --all -f -y
14
15 # Separate instructions per scope
16 RUN mkdir data-sci-demo
17
18 COPY ./your-project data-sci-demo/
```

[https://docs.docker.com/develop/develop-images/dockerfile\\_best-practices/](https://docs.docker.com/develop/develop-images/dockerfile_best-practices/)



ixek | <https://bit.ly/europython-ml-docker>

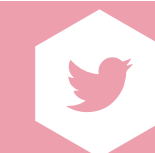


# MOUNT VOLUMES TO ACCESS DATA

- You can use bind mounts to directories (unless you are using a database)
- Avoid issues by creating a non-root user

```
docker-ds - workflow.yml
1 # mount directory
2 docker run --volume /home/user/yourproject:/yourproject mycontainer
3 # mount directory as read-only
4 docker run --volume /home/user/yourproject:/yourproject:ro mycontainer
5 # mount multiple directories, one with write access relative to current path
  (Linux)
6 docker run --volume /home/user/article-x-supplement/data:/data:ro \
7 --volume $(pwd)/outputs:/output-data:rw mycontainer
```

[https://docs.docker.com/develop/develop-images/dockerfile\\_best-practices/](https://docs.docker.com/develop/develop-images/dockerfile_best-practices/)



ixek | <https://bit.ly/europython-ml-docker>

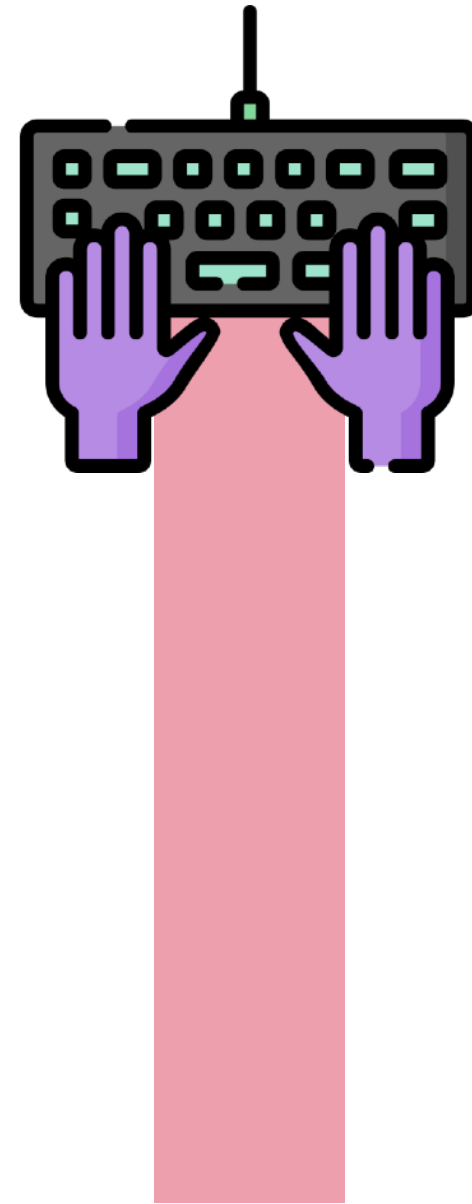


# SECURITY AND PERFORMANCE

# MINIMISE PRIVILEGE - FAVOUR LESS PRIVILEGED USER

Lock down your container:

- Run as non-root user (Docker runs as root by default)
- Minimise capabilities



```
docker-ds - Dockerfile
1 FROM python:3.8.2-slim-buster
2 RUN useradd --create-home jovyan
3 WORKDIR /home/jovyan
4 USER jovyan
```



# DON'T LEAK SENSITIVE INFORMATION

Remember Docker images are like onions. If you copy keys in an intermediate layer **they are cached**.



Keep them out of your **Dockerfile**.

# USE MULTI STAGE BUILDS

- Fetch and manage secrets in an intermediate layer
- Not all your dependencies will have been packed as wheels so you might need a compiler – build a compile and a runtime image
- Smaller images overall

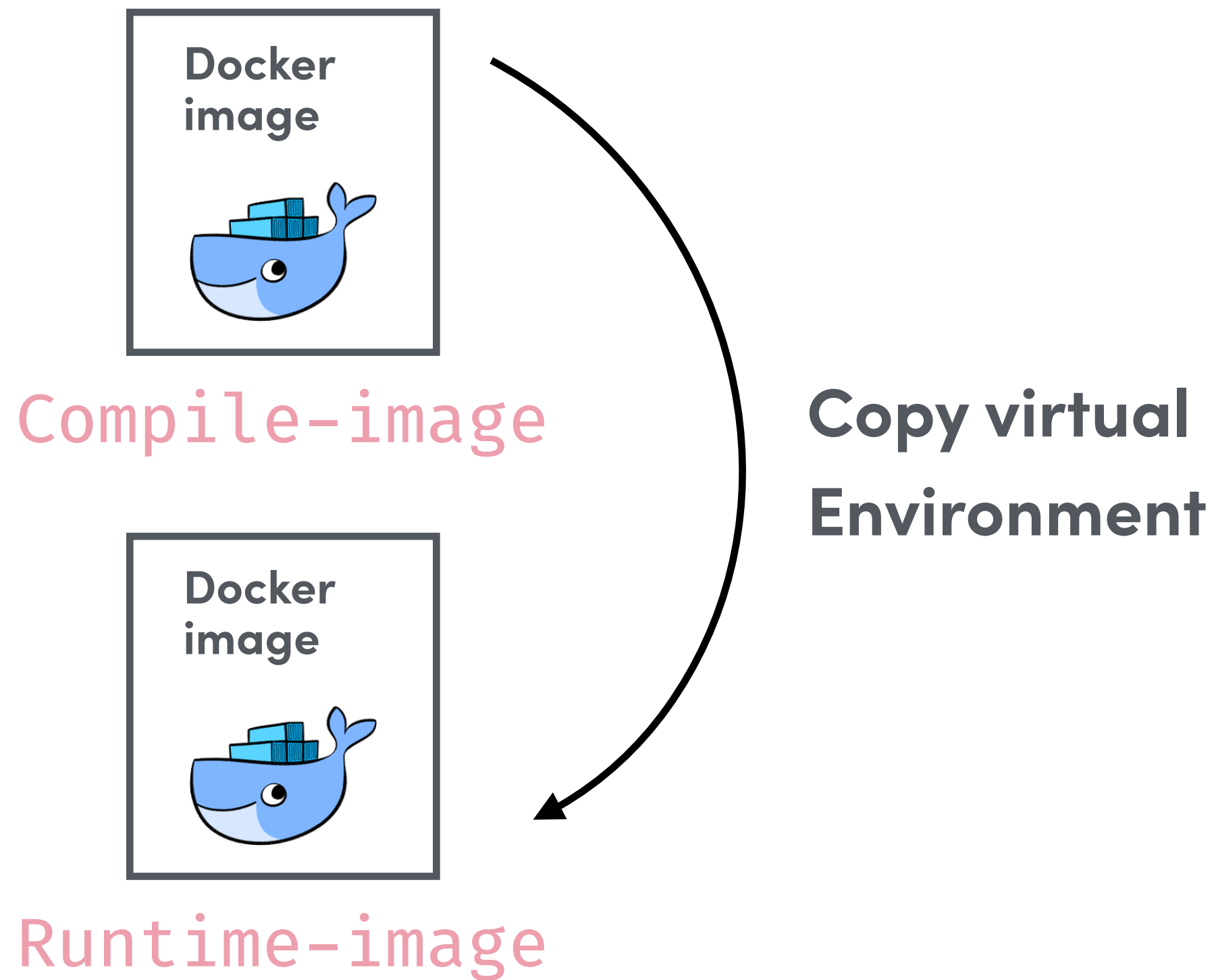


```
docker-ds - Dockerfile
1 # Always use a concrete tag (avoid LATEST)
2 FROM python:3.8.2-slim-buster as compile-image
3
4 # Add metadata
5 LABEL maintainer="Tania Allard"
6 LABEL securitytxt=
  "https://www.example.com/.well-known/security.txt"
7
8 RUN apt-get update
9 RUN
  apt-get install -y --no-install-recommends gcc build-essen
  tial gcc gfortran
10
11 RUN python -m venv /opt/venv
12
13 # Ensure we use the virtualenv
14 ENV PATH="/opt/venv/bin:$PATH"
15
16 COPY requirements.txt /tmp/
17
18 RUN CFLAGS=
  "-g0 -Wl,--strip-all -I/usr/include:/usr/local/include -L/u
  sr/lib:/usr/local/lib"
  \
19   pip install \
20   --no-cache-dir \
21   --compile \
22   --global-option=build_ext \
23   --global-option="-j 4" \
24   -r /tmp/requirements.txt
25
26 # -----
27 # This is the second image that copies the compiled librar
  y
28
29 FROM python:3.8.2-slim-buster as runtime-image
30
31 COPY --from=compile-image /opt/venv /opt/venv
32 # Ensure we use the virtualenv
33 ENV PATH="/opt/venv/bin:$PATH"
```



# USE MULTI STAGE BUILDS

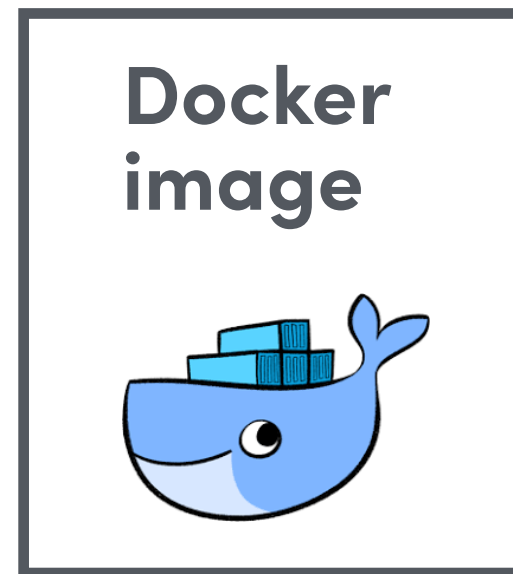
```
$ docker build --pull --rm -f "Dockerfile" \
-t trallard:data-scratch-1.0 "."
```



```
docker-ds - Dockerfile
1 # Always use a concrete tag (avoid LATEST)
2 FROM python:3.8.2-slim-buster as compile-image
3
4 # Add metadata
5 LABEL maintainer="Tania Allard"
6 LABEL securitytxt=
  "https://www.example.com/.well-known/security.txt"
7
8 RUN apt-get update
9 RUN
  apt-get install -y --no-install-recommends gcc build-essen
  tial gcc gfortran
10
11 RUN python -m venv /opt/venv
12
13 # Ensure we use the virtualenv
14 ENV PATH="/opt/venv/bin:$PATH"
15
16 COPY requirements.txt /tmp/
17
18 RUN CFLAGS=
  "-g0 -Wl,--strip-all -I/usr/include:/usr/local/include -L/u
  sr/lib:/usr/local/lib"
  \
  pip install \
  --no-cache-dir \
  --compile \
  --global-option=build_ext \
  --global-option="-j 4" \
  -r /tmp/requirements.txt
19
20
21
22
23
24
25
26 # -----
27 # This is the second image that copies the compiled librar
  y
28
29 FROM python:3.8.2-slim-buster as runtime-image
30
31 COPY --from=compile-image /opt/venv /opt/venv
32 # Ensure we use the virtualenv
33 ENV PATH="/opt/venv/bin:$PATH"
```

# USE MULTI STAGE BUILDS

## FINAL IMAGE



Runtime-image

trallard:data-scratch-1.0

```
docker-ds - Dockerfile
1 # Always use a concrete tag (avoid LATEST)
2 FROM python:3.8.2-slim-buster as compile-image
3
4 # Add metadata
5 LABEL maintainer="Tania Allard"
6 LABEL securitytxt=
  "https://www.example.com/.well-known/security.txt"
7
8 RUN apt-get update
9 RUN
  apt-get install -y --no-install-recommends gcc build-essen
  tial gcc gfortran
10
11 RUN python -m venv /opt/venv
12
13 # Ensure we use the virtualenv
14 ENV PATH="/opt/venv/bin:$PATH"
15
16 COPY requirements.txt /tmp/
17
18 RUN CFLAGS=
  "-g0 -Wl,--strip-all -I/usr/include:/usr/local/include -L/u
  sr/lib:/usr/local/lib"
  \
  pip install \
  --no-cache-dir \
  --compile \
  --global-option=build_ext \
  --global-option="-j 4" \
  -r /tmp/requirements.txt
19
20 # -----
21 # This is the second image that copies the compiled librar
  y
22
23 FROM python:3.8.2-slim-buster as runtime-image
24
25 COPY --from=compile-image /opt/venv /opt/venv
26 # Ensure we use the virtualenv
27 ENV PATH="/opt/venv/bin:$PATH"
```



**AUTOMATE**



# PROJECT TEMPLATES



Need a standard project template?

Use `cookie cutter data science`

Or `cookie cutter docker science`

```
docker-ds - workflow.yml
1  |— Makefile           ← Makefile contains many targets such as create docker container or
2  |                   get input files.
3  |— config           ← This directory contains configuration files used in scripts
4  |   |— jupyter_config.py or Jupyter Notebook.
5  |   |
6  |   |— data         ← data directory contains the input resources.
7  |   |— docker      ← docker directory contains Dockerfile.
8  |   |   |— Dockerfile ← Dockerfile have the container settings. Users modify Dockerfile
9  |   |   |           if additional library is needed for experiments.
10 |   |— model        ← model directory store the model files created in the experiments.
11 |   |— my_data_science_project ← cookie-cutter-docker-science creates the directory whose name is same
12 |   |   |           as project name. In this directory users puts python files used in scripts
13 |   |   |           or Jupyter Notebook.
14 |   |   |— __init__.py
15 |   |— notebook    ← This directory stores the ipynb files saved in Jupyter Notebook.
16 |   |— requirements.txt ← Libraries needed to run experiments. The library listed in this file
17 |   |               are installed in the Docker container.
18 |   |— scripts      ← Users add the script files to generate model files or run evaluation.
```

<https://github.com/docker-science/cookiecutter-docker-science>

<https://drivendata.github.io/cookiecutter-data-science/>

```
$ conda install jupyter repo2docker
$ jupyter-repo2docker “.”
```



# DO NOT REINVENT THE WHEEL

Leverage the existence and usage of tools like `repo2docker`.

Already configured and optimised for Data Science / Scientific computing.

# DO NOT REINVENT THE WHEEL

Leverage the existence and usage of tools like `repo2docker`.

Already configured and optimised for Data Science / Scientific computing.

- Configuration Files
  - `environment.yml` - Install a Python environment
  - `Pipfile` and/or `Pipfile.lock` - Install a Python environment
  - `requirements.txt` - Install a Python environment
  - `setup.py` - Install Python packages
  - `Project.toml` - Install a Julia environment
  - `REQUIRE` - Install a Julia environment (legacy)
  - `install.R` - Install an R/RStudio environment
  - `apt.txt` - Install packages with apt-get
  - `DESCRIPTION` - Install an R package
  - `manifest.xml` - Install Stencila
  - `postBuild` - Run code after installing the environment
  - `start` - Run code before the user sessions starts
  - `runtime.txt` - Specifying runtimes
  - `default.nix` - the nix package manager
  - `Dockerfile` - Advanced environments

# DELEGATE TO YOUR CONTINUOUS INTEGRATION TOOL

Set Continuous integration  
(Travis, GitHub Actions, whatever  
you prefer).

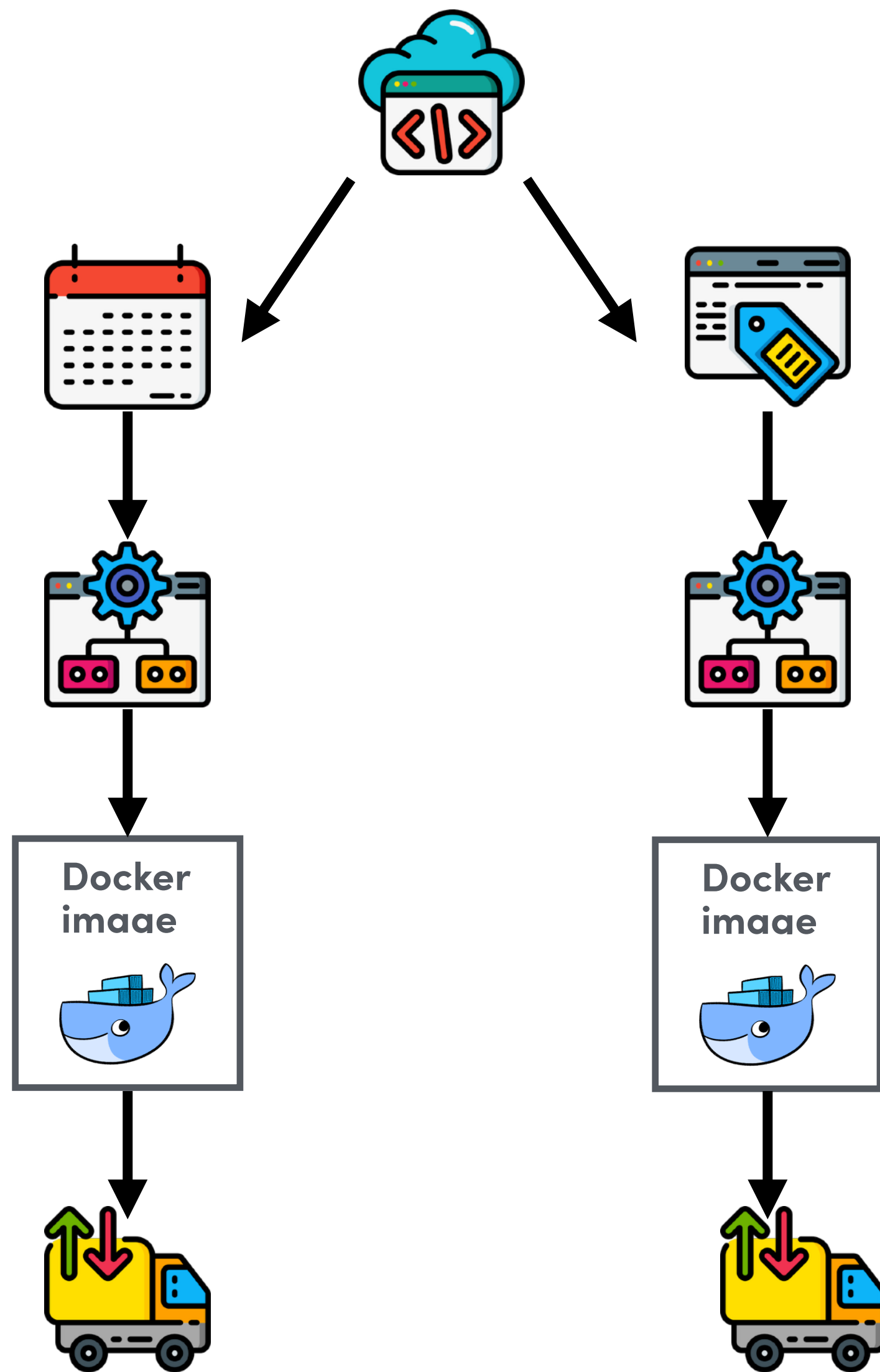
And delegate your build - also  
build often.

docker-ds - workflow.yml

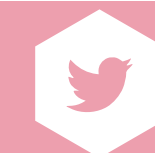
```
1 name: Publish to Registry
2 on:
3   release:
4     types: [published]
5   schedule:
6     # Build your images frequently
7     - cron: "0 2 * * 0" # Weekly on Sundays at 02:00
8 jobs:
9   update:
10    runs-on: ubuntu-latest
11    steps:
12      - uses: actions/checkout@master
13      - name: Get release version
14        id: get_version
15        run: echo ::set-env name=RELEASE_VERSION::$(echo ${GITHUB_REF:10})
16      - name: Build and publish
17        uses: docker/build-push-action@v1
18        with:
19          username: ${ secrets.DOCKER_USERNAME }
20          password: ${ secrets.DOCKER_PASSWORD }
21          repository: myorg/myrepository
22          tag_with_ref: true
23          tag_with_sha: true
24
```



# THIS WORKFLOW



- Code in version control
- Trigger on tag / Also scheduled trigger
- Build image
- Push image





# TOP TIPS

# TOP TIPS

1. **Rebuild your images frequently - get security updates for system packages**
2. **Never work as root / minimise the privileges**
3. **You do not want to use Alpine Linux (go for buster, stretch or the Jupyter stack)**
4. **Always know what you are expecting: pin / version EVERYTHING (use pip-tools, conda, poetry or pipenv)**
5. **Leverage build cache**

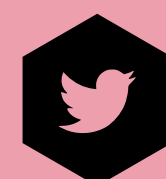
# TOP TIPS

6. Use one Dockerfile per project
7. Use multi-stage builds - need to compile code? Need to reduce your image size?
8. Make your images identifiable (test, production, R&D) - also be careful when accessing databases and using ENV variables / build variables
9. Do not reinvent the wheel! Use repo2docker
10. Automate - no need to build and push manually
11. Use a linter





# THANK YOU



@ixek



@trallard



trallard.dev